

Multi-Classification by Categorical Features via Clustering

Yevgeny Seldin

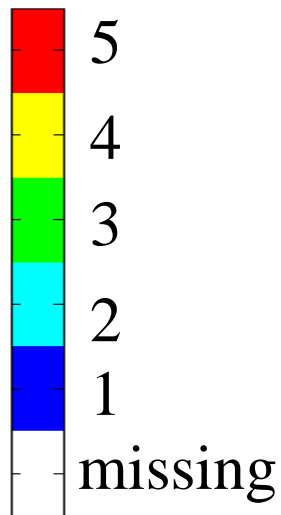
Joint work with Naftali Tishby

The Hebrew University of Jerusalem

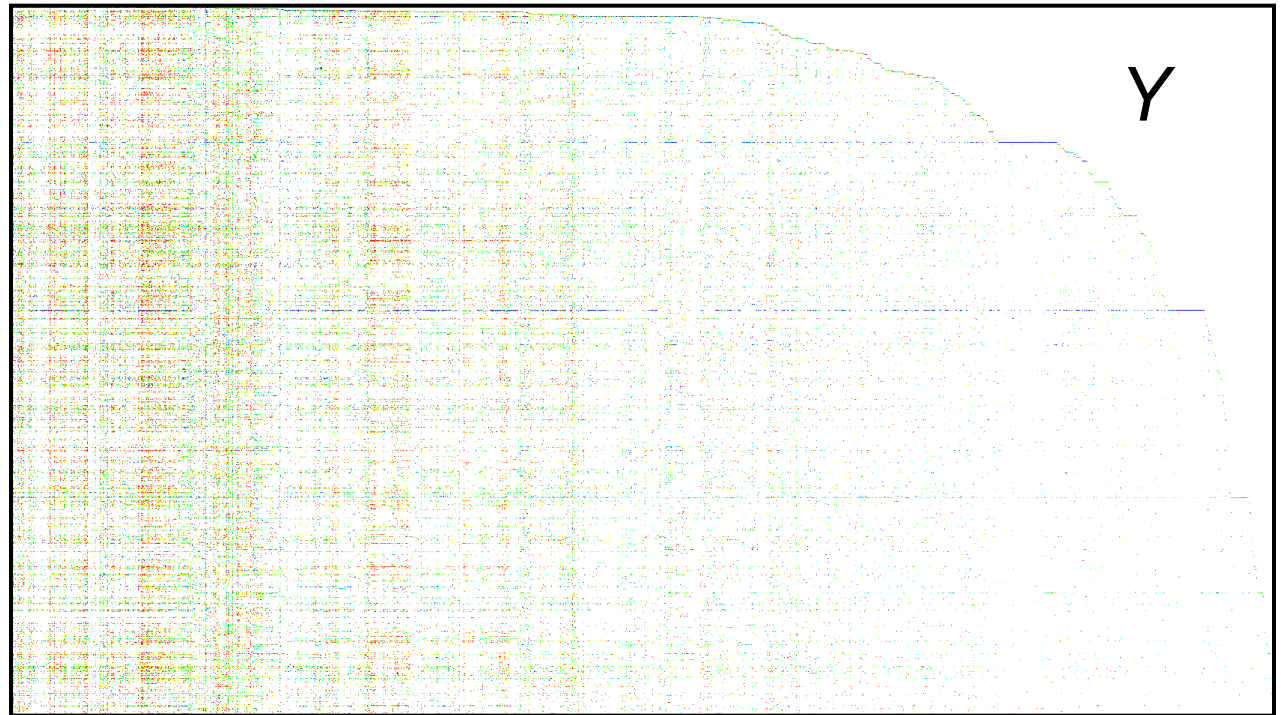
Multi-Classification by Categorical Features

Example: Collaborative Filtering

Ratings
bar



Viewers (X_1)

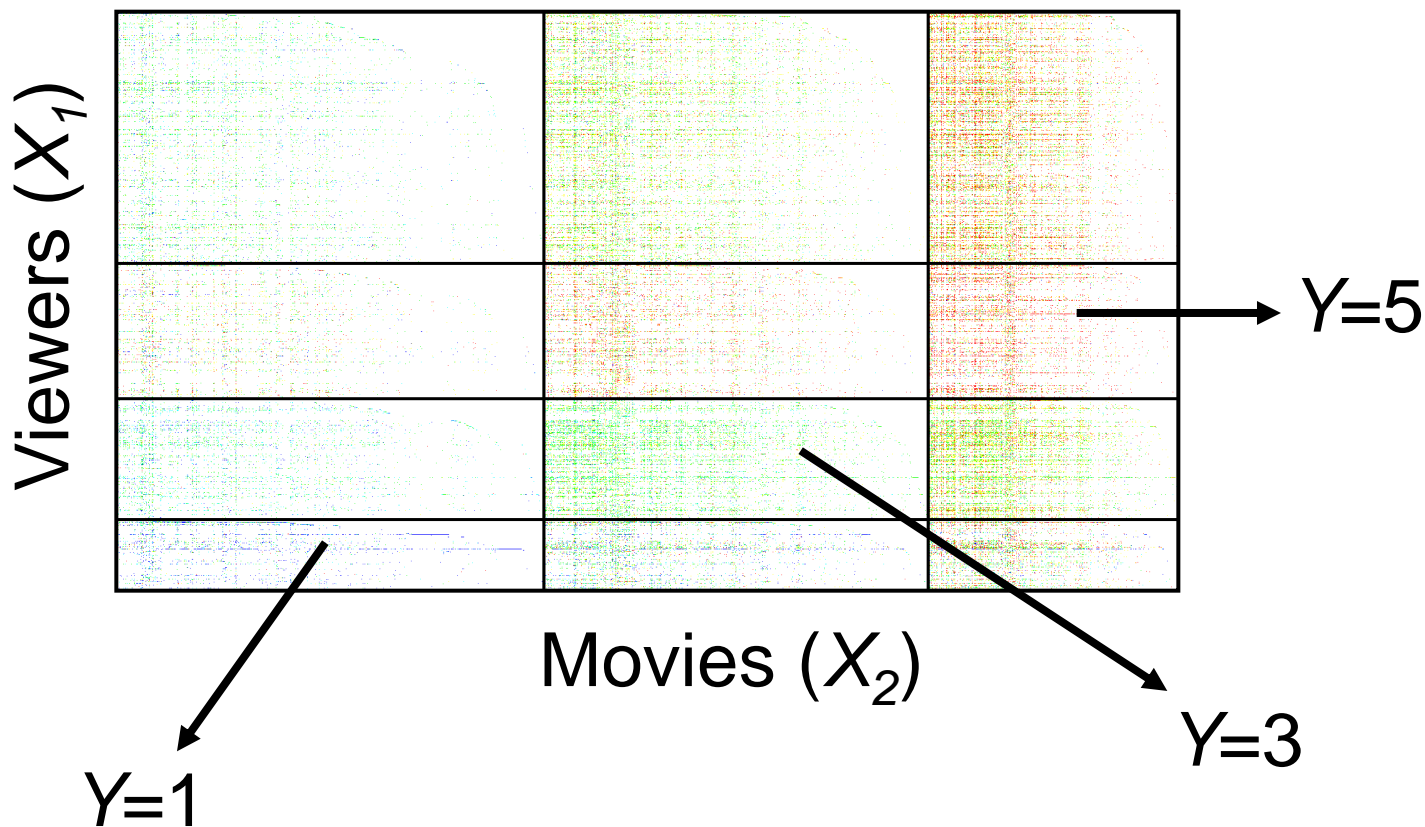
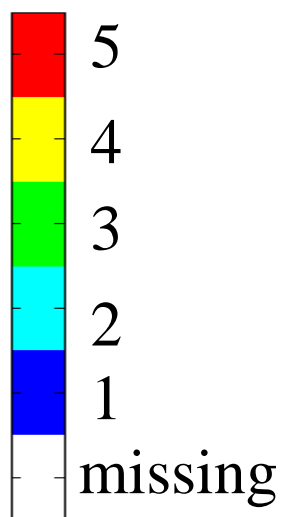


MovieLens Data

Movies (X_2)

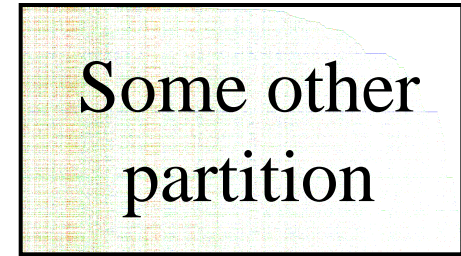
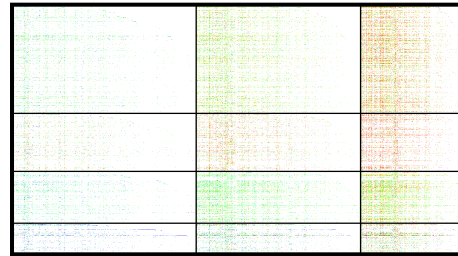
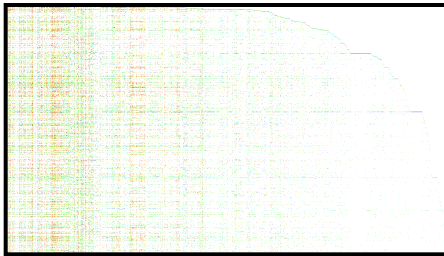
Multi-Classification via Grid Clustering

Ratings
bar



[Seldin, Slonim & Tishby, NIPS06]

Question: which partition is better?



- Statistical Reliability vs. Precision Tradeoff

The Approach We Take

- Relate with generalization properties

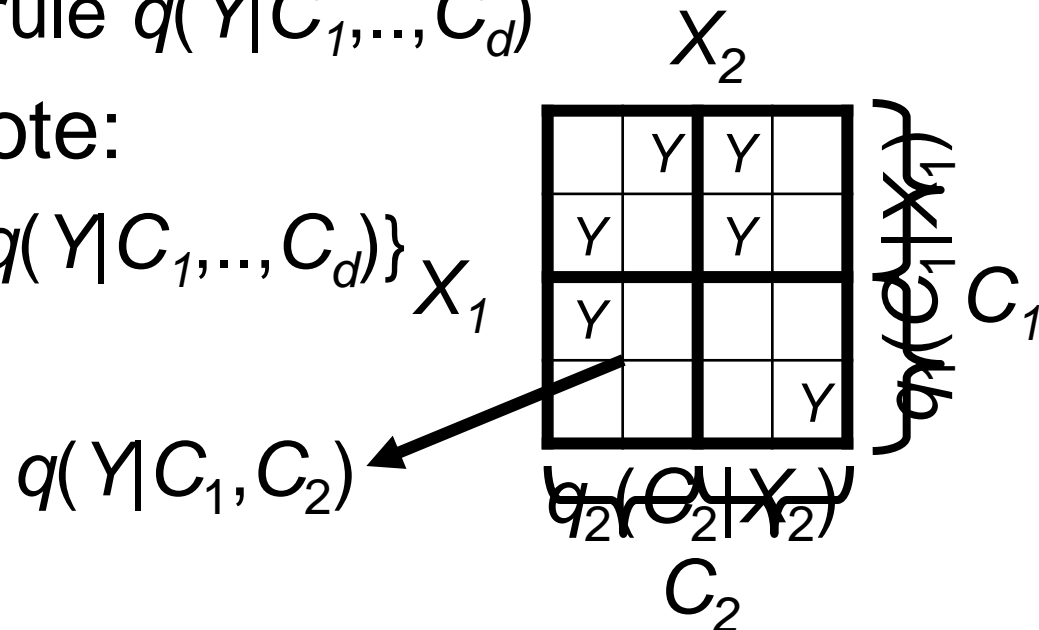
Some Definitions

- Classification via Stochastic Grid Clustering:

- A set of stochastic mappings $q_i(C_i|X_i)$
- A classification rule $q(Y|C_1, \dots, C_d)$

- Collectively denote:

- $Q = \{\{q_i(C_i|X_i)\}, q(Y|C_1, \dots, C_d)\}$



Generalization Bound

- With probability $\geq 1-\delta$:

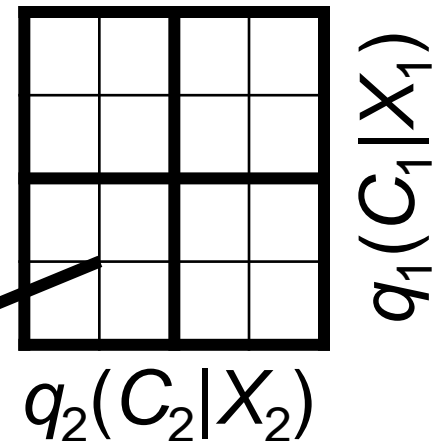
$$L(Q) \leq \hat{L}(Q) + \sqrt{\frac{\sum_i n_i I_U(X_i; C_i) + K}{2N}}$$

Expected
loss of Q

Empirical
loss of Q

Complexity
term

$q(Y|C_1, C_2)$

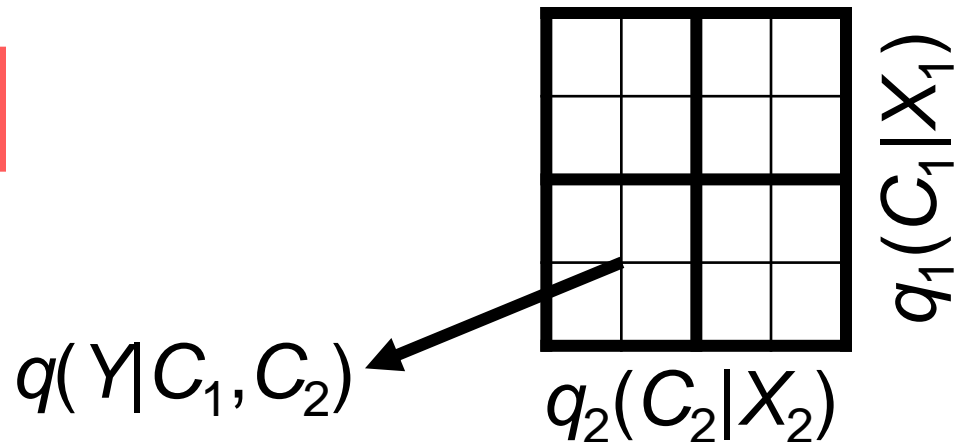


Generalization Bound

- With probability $\geq 1-\delta$:

$$L(Q) \leq \hat{L}(Q) + \sqrt{\frac{\sum_i n_i I_U(X_i; C_i) + K}{2N}}$$

$$n_i = |X_i|$$



Generalization Bound

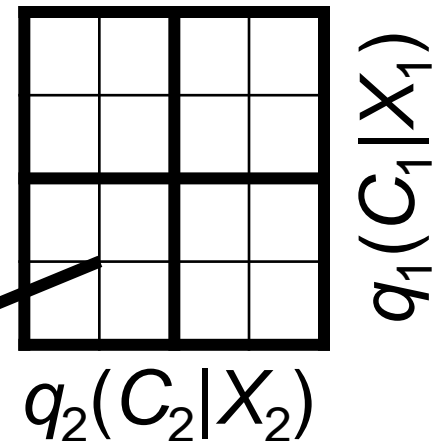
- With probability $\geq 1-\delta$:

$$L(Q) \leq \hat{L}(Q) + \sqrt{\frac{\sum_i n_i I_U(X_i; C_i) + K}{2N}}$$

Information
preserved by
 $q_i(C_i|X_i)$

(w.r.t. $p(X_i)=1/n_i$)

$q(Y|C_1, C_2)$

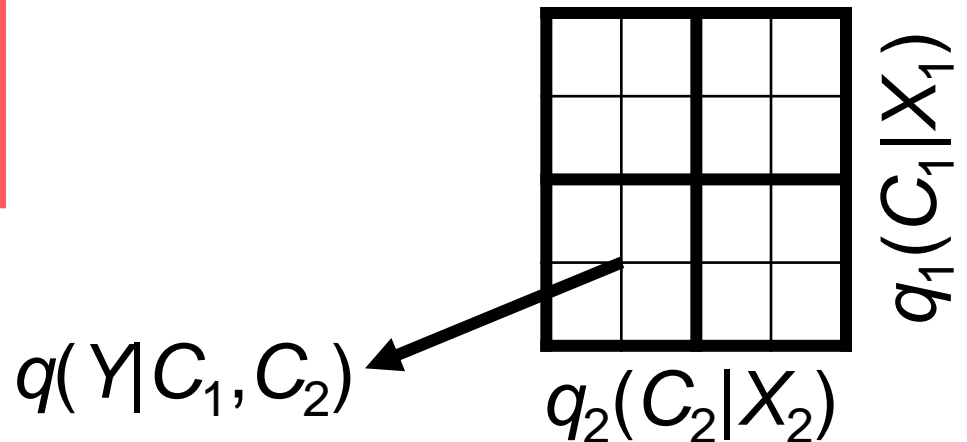


Generalization Bound

- With probability $\geq 1-\delta$:

$$L(Q) \leq \hat{L}(Q) + \sqrt{\frac{\sum_i n_i I_U(X_i; C_i) + K}{2N}}$$

Independent of
 $q_i(C_i|X_i)$

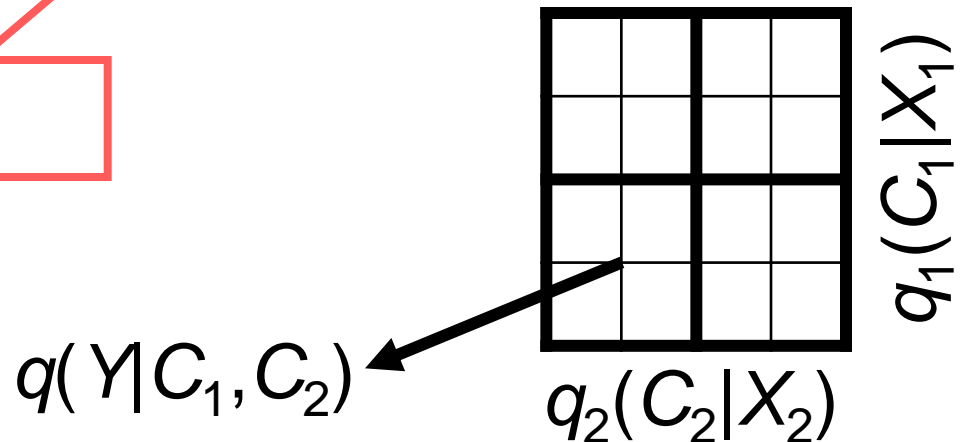


Generalization Bound

- With probability $\geq 1-\delta$:

$$L(Q) \leq \hat{L}(Q) + \sqrt{\frac{\sum_i n_i I_U(X_i; C_i) + K}{2N}}$$

Sample size

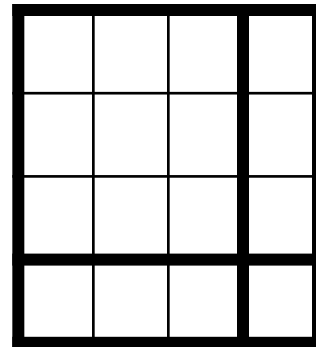


Generalization Bound

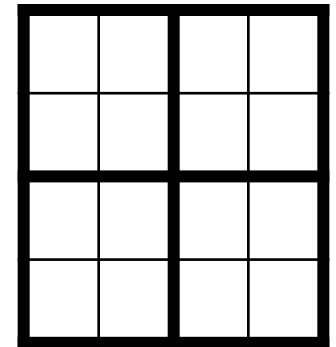
- With probability $\geq 1-\delta$:

$$L(Q) \leq \hat{L}(Q) + \sqrt{\frac{\sum_i n_i I_U(X_i; C_i) + K}{2N}}$$

Optimization tradeoff:
Empirical loss vs.
“Effective” partition
complexity



Lower
Complexity



Higher
Complexity

Generalization Bound

- With probability $\geq 1-\delta$:

$$L(Q) \leq \hat{L}(Q) + \sqrt{\frac{\sum_i n_i I_U(X_i; C_i) + K}{2N}}$$

$$K = \underbrace{\sum_i \left(m_i \ln n_i + \frac{(\ln n_i + 1)^2}{4} \right)}_{\text{Logarithmic in } n_i} + \underbrace{\left(\prod_i m_i \right)}_{\text{Number of partition cells}} \ln |Y| + \underbrace{\frac{1}{2} \ln(4N) + \ln \frac{1}{\delta}}_{\text{"usual stuff"}}$$

$$m_i = |C_i|$$

Logarithmic
in n_i

Number of
partition cells

"usual stuff"

Proof Idea

- Start with the PAC-Bayesian bound:

$$L(Q) \leq \hat{L}(Q) + \sqrt{\frac{D(Q \| P) + \ln(4N)/2 + \ln(1/\delta)}{2N}}$$

– [McAllester 99], [Maurer 04]

Proof Idea

- Start with the PAC-Bayesian bound:

$$L(Q) \leq \hat{L}(Q) + \sqrt{\frac{D(Q \| P) + \ln(4N)/2 + \ln(1/\delta)}{2N}}$$

– [McAllester 99], [Maurer 04]

- Design a combinatorial prior $P(h)$ by counting the number of hard partitions

Proof Idea

- Start with the PAC-Bayesian bound:

$$L(Q) \leq \hat{L}(Q) + \sqrt{\frac{D(Q \| P) + \ln(4N)/2 + \ln(1/\delta)}{2N}}$$

– [McAllester 99], [Maurer 04]

- Design a combinatorial prior $P(h)$ by counting the number of hard partitions
- Calculate $D(Q \| P)$

Proof Idea

- Start with the PAC-Bayesian bound:

$$L(Q) \leq \hat{L}(Q) + \sqrt{\frac{D(Q \| P) + \ln(4N)/2 + \ln(1/\delta)}{2N}}$$

– [McAllester 99], [Maurer 04]

- Design a combinatorial prior $P(h)$ by counting the number of hard partitions
- Calculate $D(Q \| P)$
- Details: at the paper/poster

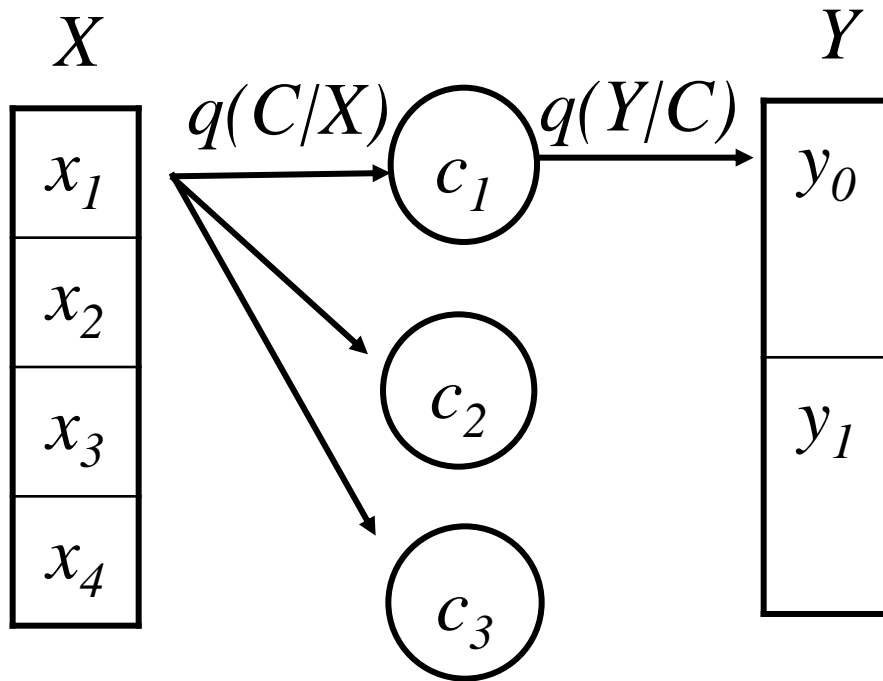
Messages

- For Clustering:
 - Evaluate clustering by its generalization properties on the task it is designed for
- For Classification:
 - Unify feature values to amplify statistical reliability

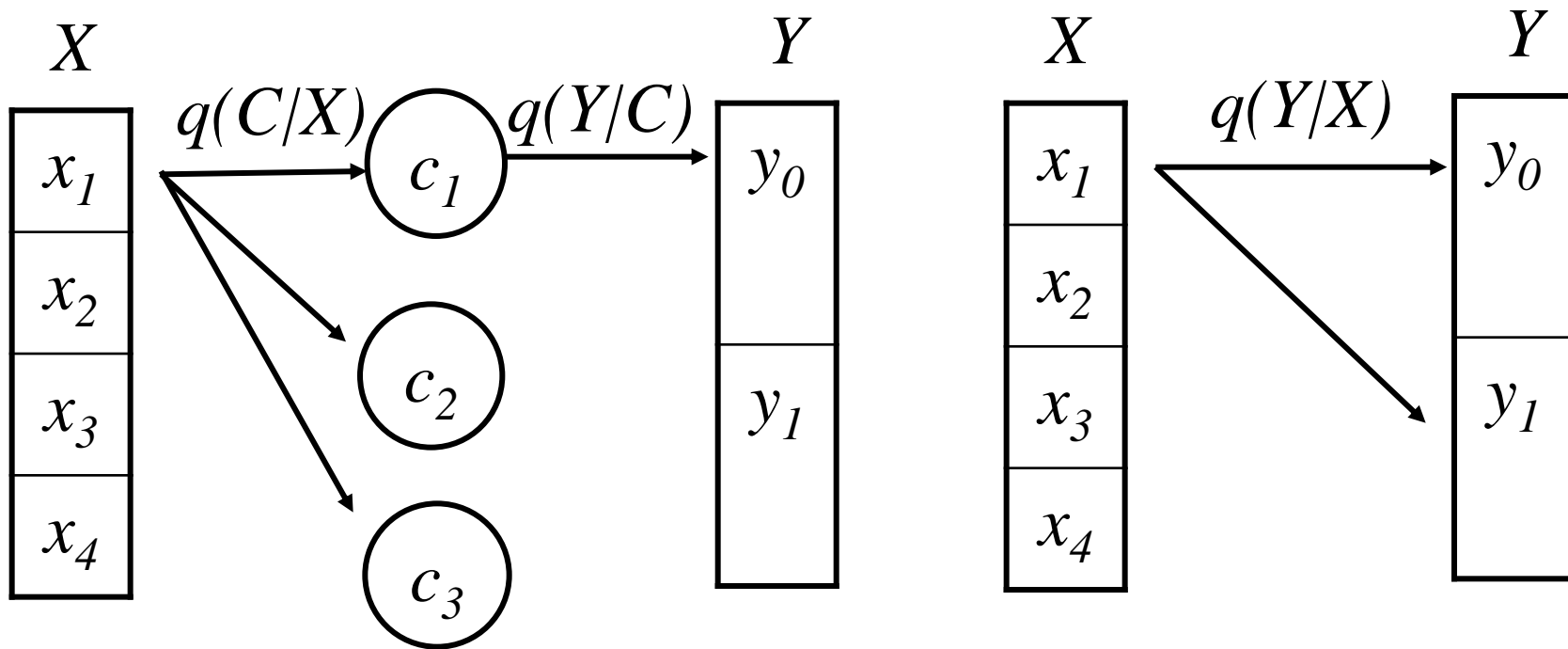
Classification by a Single Feature

- A tighter and simpler bound
- Application: Feature Ranking

Classification by a Single Feature

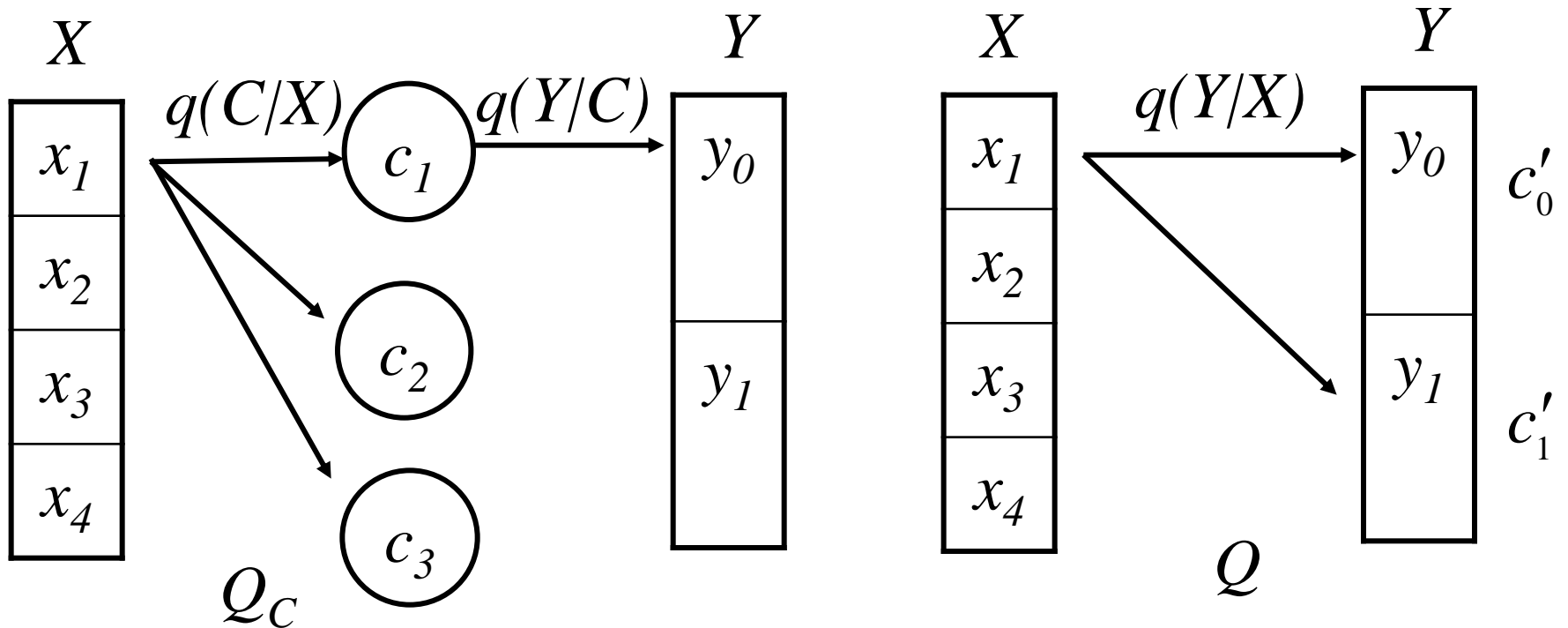


Equivalent “Direct Mapping”



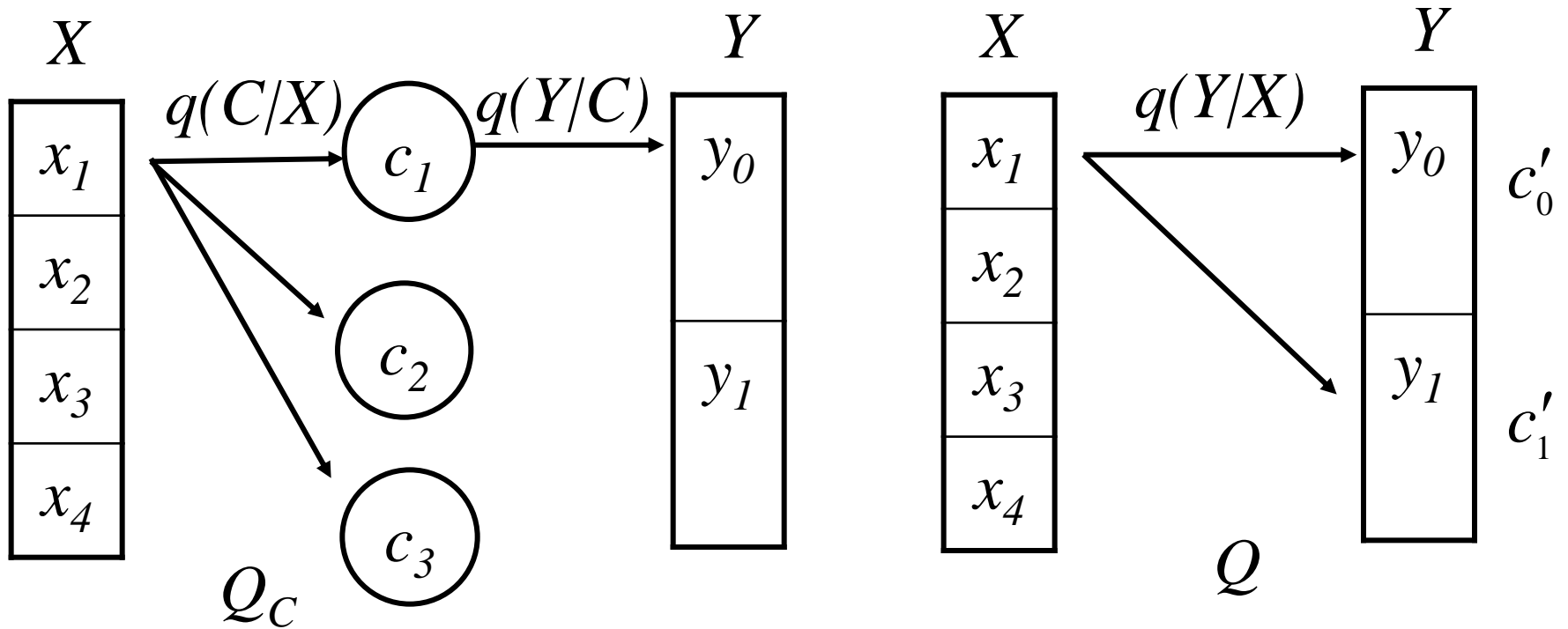
Define : $q(y/x) = \sum_c q(c|x)q(y|c)$

Optimality of Direct Mappings



Define : $q(y/x) = \sum_c q(c|x)q(y|c)$

Optimality of Direct Mappings



$$\text{Define : } q(y/x) = \sum_c q(c|x)q(y|c)$$

$$\hat{L}(Q_C) = \hat{L}(Q), \quad I_U(X;C) \geq I_U(X;Y)$$

A Bound for Direct Mappings

$$L(Q) \leq \hat{L}(Q) + \sqrt{\frac{nI_U(X;Y) + K'}{2N}}$$

A Bound for Direct Mappings

$$L(Q) \leq \hat{L}(Q) + \sqrt{\frac{nI_U(X;Y) + K'}{2N}}$$

- Tighter than the bound on $L(Q_C)$

A Bound for Direct Mappings

$$L(Q) \leq \hat{L}(Q) + \sqrt{\frac{nI_U(X;Y) + K'}{2N}}$$

- Tighter than the bound on $L(Q_C)$
- Holds for any classification rule $q(Y|X)$

A Bound for Direct Mappings

$$L(Q) \leq \hat{L}(Q) + \sqrt{\frac{nI_U(X;Y) + K'}{2N}}$$

- Tighter than the bound on $L(Q_C)$
- Holds for any classification rule $q(Y|X)$
- Can be optimized (gradient descent) with respect to $q(Y|X)$ to provide an optimal classification rule

A Bound for Direct Mappings

$$L(Q) \leq \hat{L}(Q) + \sqrt{\frac{nI_U(X;Y) + K'}{2N}}$$

- Tighter than the bound on $L(Q_C)$
- Holds for any classification rule $q(Y|X)$
- Can be optimized (gradient descent) with respect to $q(Y|X)$ to provide an optimal classification rule
- No need for intermediate clustering
 - Only for a single feature

Some Insights

$$L(Q) \leq \hat{L}(Q) + \sqrt{\frac{nI_U(X;Y) + K'}{2N}}$$

- $\hat{L}(Q)$ is minimized by $q_{ml}(y|x)$

Some Insights

$$L(Q) \leq \hat{L}(Q) + \sqrt{\frac{nI_U(X;Y) + K'}{2N}}$$

- $\hat{L}(Q)$ is minimized by $q_{ml}(y|x)$
- For $I_U(X;Y)=0$, $L(Q)$ is minimized by $q_{ml}(y)$

Some Insights

$$L(Q) \leq \hat{L}(Q) + \sqrt{\frac{nI_U(X;Y) + K'}{2N}}$$

- $\hat{L}(Q)$ is minimized by $q_{ml}(y|x)$
- For $I_U(X;Y)=0$, $L(Q)$ is minimized by $q_{ml}(y)$
- Thus $L(Q)$ is minimized by smoothing $q_{ml}(y|x)$ toward $q_{ml}(y)$

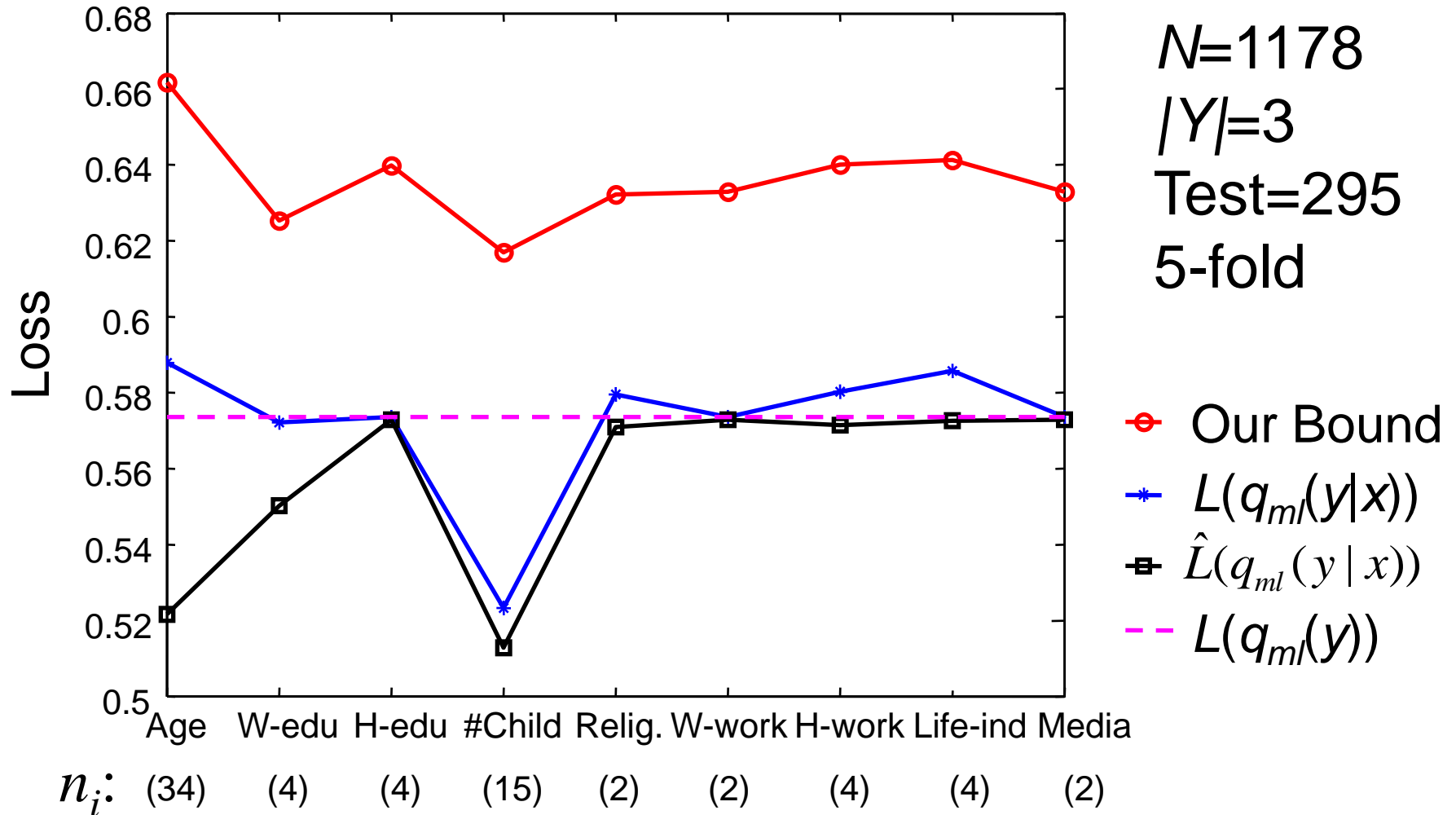
Application: Feature Ranking

- Rank features by their generalization potential
 - And not mutual information or correlation with the label
- Especially important for features of different cardinalities and small sample
 - Example: $Y = \text{cancer/no_cancer}$
 $X_1 = \text{smoking/not_smoking}$
 $X_2 = \text{year_of_birth}$

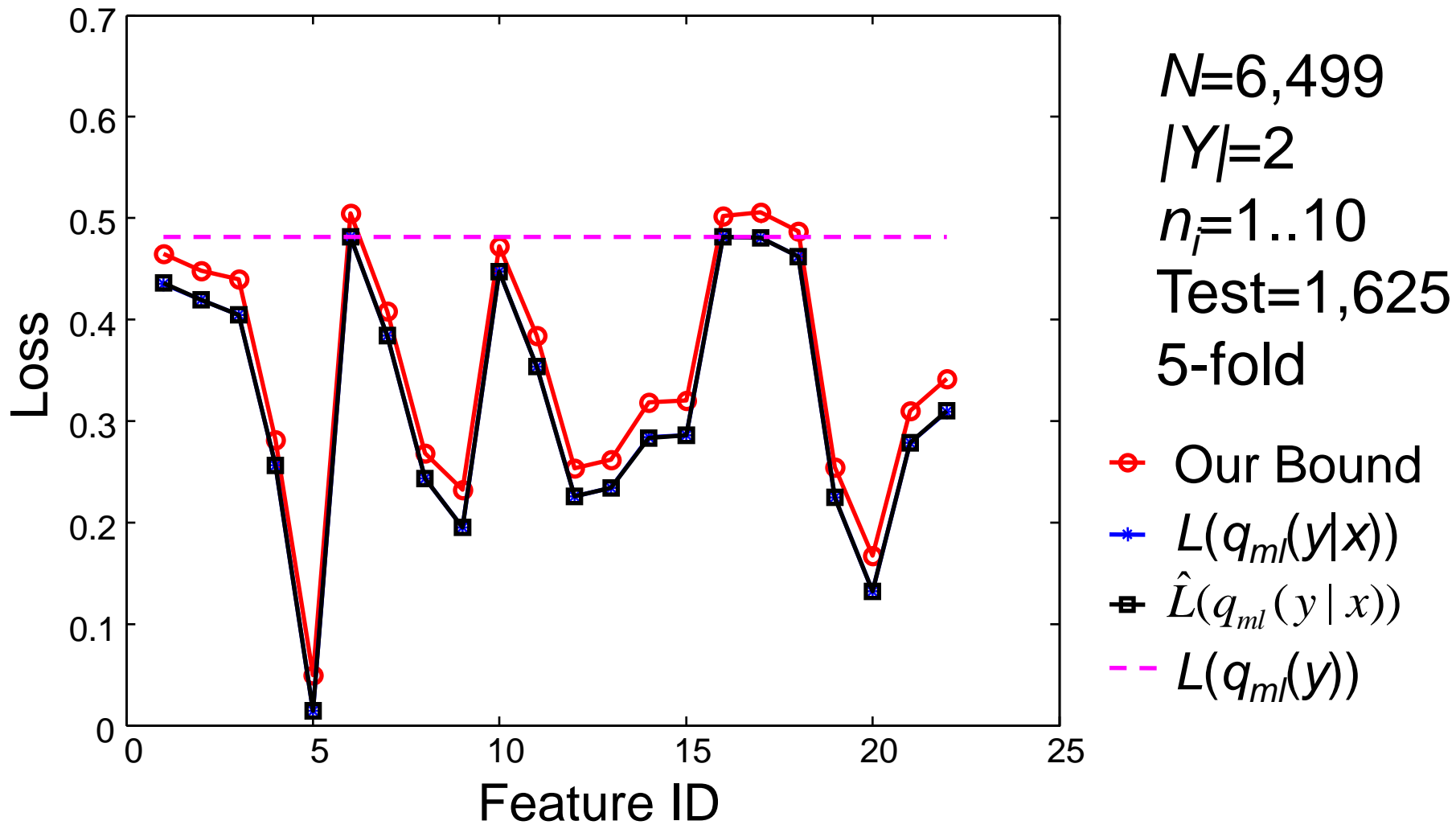
Related work

- Sabato, Shalev-Shwartz, COLT07
 - Generalization bound for $q(Y|X)=\hat{p}(Y|X)$
(empirical distribution)
- Our work:
 - Any $q(Y|X)$, in particular $q_{ml}(Y|X)$

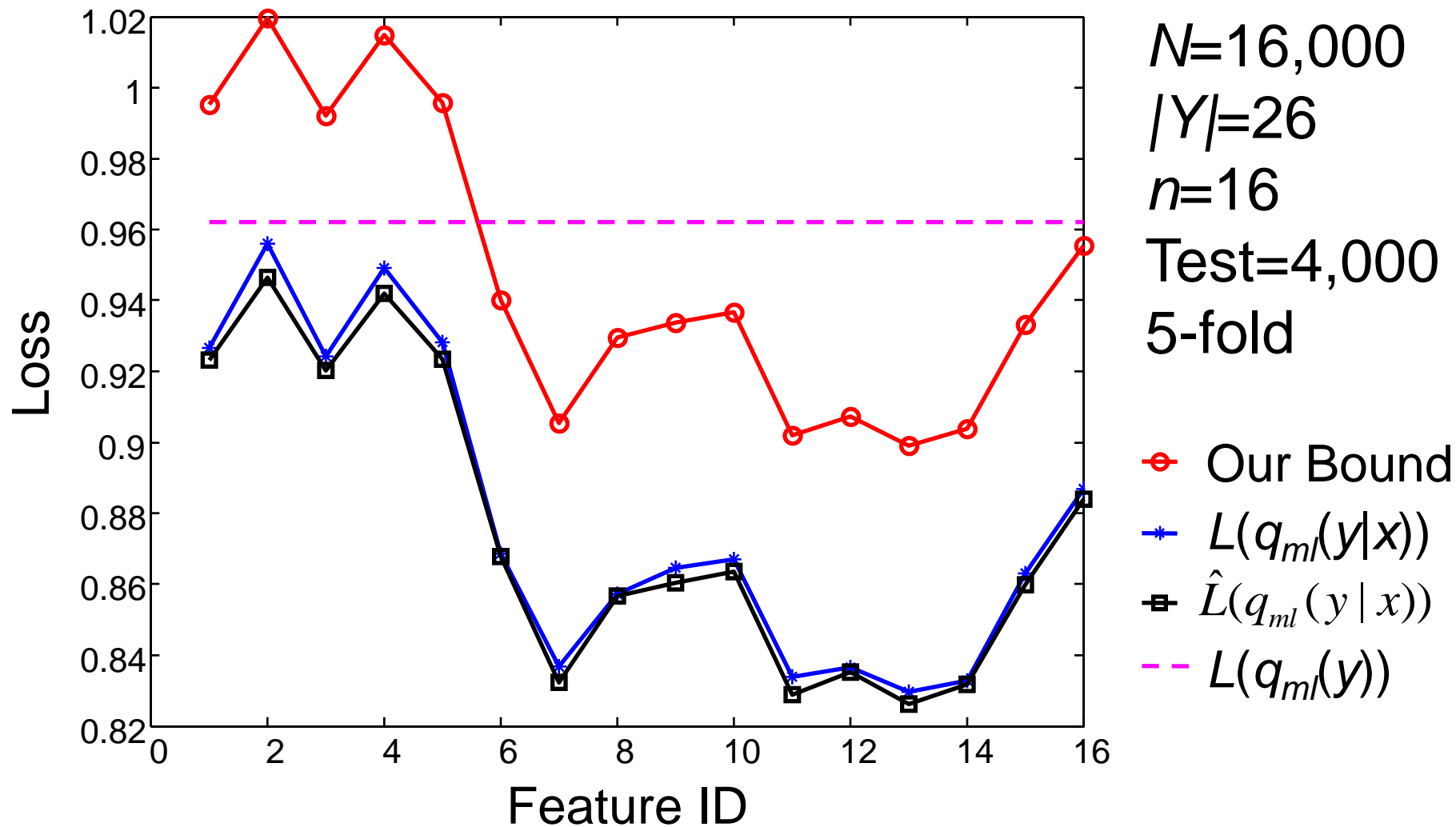
Contraceptive Method Choice



Mushrooms



Letters

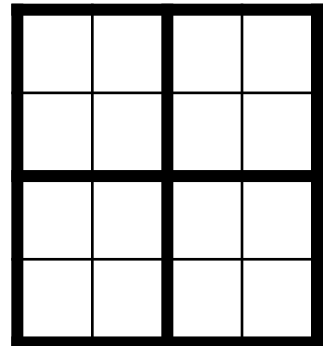


Comparison with MI and Corr



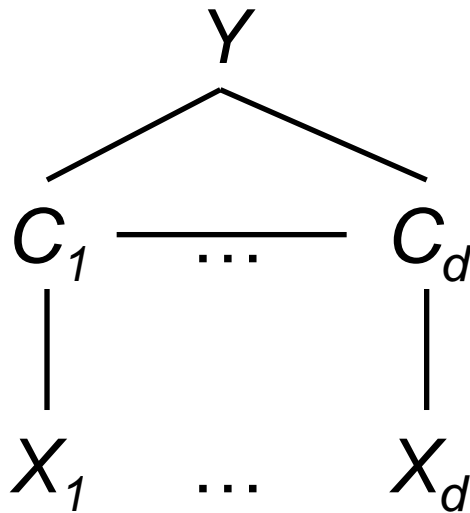
Summary

- Generalization bound for multi-classification based on grid clustering
- Unify feature values to amplify statistics
- Evaluation of clustering by its generalization power on a given task
- Feature Ranking
- Limitation: high dimensions

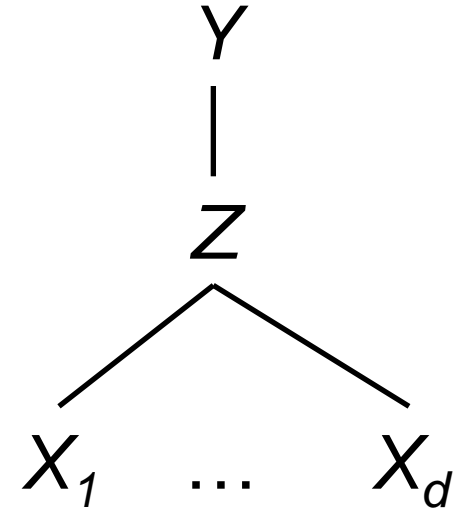
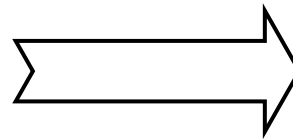


Future Work

- Derive a generalization bound for general graphical models



Grid Clustering



Factor Model